

SYSTEM AND METHOD FOR SERVING A WEB
SITE FROM MULTIPLE SERVERS

BACKGROUND OF THE INVENTION

1. Related Applications.

- 5 The present invention claims priority from U.S. Provisional Patent Application No. 60/197,490 entitled CONDUCTOR GATEWAY filed on April 17, 2000.

2. Field of the Invention.

- 10 The present invention relates, in general, to network information access and, more particularly, to software, systems and methods for serving web pages in a coordinated fashion from multiple cooperating web servers.

3. Relevant Background.

- 15 Increasingly, business data processing systems, entertainment systems, and personal communications systems are implemented by computers across networks that are interconnected by Internetworks (e.g., the Internet). The Internet is rapidly emerging as the preferred system for distributing and exchanging data. Data exchanges support
20 applications including e-commerce, broadcast and multicast messaging, videoconferencing, gaming, and the like.

- The Internet is a collection of disparate computers and networks coupled together by a web of interconnections using standardized communications protocols. The Internet
25 is characterized by its vast reach as a result of its wide and increasing availability and easy access protocols.

Unfortunately, the heterogenous nature of the Internet results in variable bandwidth and quality of service between points. The latency and reliability of data transport is largely determined by the total amount of traffic on the Internet and so varies wildly seasonally and throughout the day. Other factors that affect quality of service include equipment outages and line degradation that force packets to be rerouted, damaged and/or dropped. Also, routing software and hardware limitations within the Internet infrastructure may create bandwidth bottlenecks even when the mechanisms are operating within specifications.

Internet transport protocols do not effectively discriminate between users. Data packets are passed between routers and switches that make up the Internet fabric based on the hardware's instantaneous view of the best path between source and destination nodes specified in the packet. Because each packet may take a different path, the latency of a packet cannot be guaranteed and, in practice, varies significantly. Although Internet infrastructure components have been developed that provide some level of packet prioritization, these components have limited deployment. Because of this limited deployment, web site operators cannot depend on prioritization services to be available throughout the network, and in reality can only depend on a small portion of the Internet being served by such components. Moreover, components that are available may require modifications to the web site to operate. Hence, in large part data packets are routed through the Internet without any prioritization based on content.

Prioritization has not been a major issue with conventional networks such as local area networks (LANs)

and wide area networks (WANs) that primarily handle well characterized data traffic between a known and controlled set of network nodes. In such environments, the network administrators are in control of network traffic and can effectively police bandwidth consumption, resource availability and distribution, and other variables that affect quality of service. In contrast, public networks such as the Internet, involve data traffic that is out of control of any particular network administrator. However, because of their wide availability, there is an increasing desire to provide network services using public networks, such as the Internet. At the same time, there is an increasing demand for network applications that cannot tolerate high and variable latency. This situation is complicated when the application is to be run over the Internet where latency and variability in latency are many times greater than in LAN and WAN environments.

A particular need exists in environments that involve multiple users accessing a network resource such as a web server. Examples include broadcast, multicast and videoconferencing as well as most electronic commerce (e-commerce) applications. In these applications, it is important to maintain a reliable connection so that the server and clients remain synchronized and information is not lost.

In e-commerce applications, it is important to provide a satisfying buyer experience that leads to a purchase transaction. To provide this high level of service, a web site operator must ensure that data is delivered to the customer in the most usable and efficient fashion. Also, the web site operator must ensure that critical data received from the customer is handled with priority.

Because the few techniques that provide traffic prioritization are not widely distributed in public networks, the e-commerce site owner has had little or no control over the transport mechanisms through the Internet that affect the latency and quality of service. This is akin to a retailer being forced to deal with a customer by shouting across the street, never certain how often what was said must be repeated, and knowing that during rush hour communication would be nearly impossible. While efforts are continually being made to increase the capacity and quality of service afforded by the Internet, it is contemplated that congestion will always impact the ability to predictably and reliably offer a specified level of service. Moreover, the change in the demand for bandwidth increases at a greater rate than does the change in bandwidth supply, ensuring that congestion will continue to be an issue into the foreseeable future. A need exists for a system to exchange data over the Internet that provides a high quality of service even during periods of congestion.

Many e-commerce transactions are abandoned by the user because system performance degradation frustrates the purchaser before the transaction is consummated. While a data exchange that is abandoned while a customer is merely browsing through a catalog may be tolerable, abandonment when the customer is just a few clicks away from a purchase is highly undesirable. However, existing Internet transport mechanisms and systems do not allow the e-commerce site owner any effective ability to offer differentiated service levels to, the "just browsing" and the "about-to-buy" customers. In fact, the vagaries of the Internet may lead to the casual browser receiving a higher quality of service while the about to buy customer becomes frustrated and abandons the transaction.

Partial solutions have been implemented by systems that cache Internet content at multiple geographically distributed locations. In theory, when content can be served from a cache location, it can be delivered with lower latency than if it were served from a single originating web server. However, the content must be copied from its origin to the multiple caches resulting in a tremendous volume of data that must be duplicated and transported. Moreover, it is difficult to keep all of the cache copies coherent with the origin. Furthermore, a cache lacks the intelligence to provide dynamically generated content and so is only a partial solution to the many web sites that use dynamically generated web pages.

SUMMARY OF THE INVENTION

Briefly stated, the present invention involves a system for serving web pages to a client in response to a client request specifying a resource. An originating web server having a first IP address upon which the requested resource resides is coupled to obtain content from content source. The originating web server has a request interface for receiving requests from a network and a response interface for sending responses to the requests. A communication network is coupled to the originating web server. A front-end server is provided having a second IP address, a first interface for communicating with a client application and a second interface for communicating with the originating web server. The front-end server is coupled to both obtain the request-specified resources and augment the request-specified resources using unspecified resources to generate a response to the client request from the originating web server.

Another aspect of the present invention involves a method for serving web pages to a client in response to a client request specifying a network resource. An originating web server is provided having a first IP address upon which the requested resource resides and is coupled to obtain content from content source. The originating web server has a request interface for receiving requests from a network and a response interface for sending responses to the requests. A front-end server is provided having a second IP address. A client request addressed to the originating web server is redirected to the front-end server. The front-end server is used to obtain the request-specified resources. The request-specified resources are augmented using resources within the front-end server to generate a response to the client request.

Still another aspect of the present invention involves web site delivery over a communication network. An originating web server coupled to the network executes software to access a first set of content and functionality in response to requests from a web client. A front-end server coupled to the network executes software to access a second set of content and functionality in response to requests from the web client. A first communication channel within the network supports request and response communication between the web client and the originating web server. A second communication channel within the network supports communication between the front-end server and the originating web server. Means operating cooperatively between the originating web server and the front-end server receives requests for web pages from the web site from the web client and serves web pages to the requesting web client in response to the received requests to implement the web site.

BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 illustrates a general distributed computing environment in which the present invention is implemented;

FIG. 2 shows in block-diagram form significant components of a system in accordance with the present invention;

FIG. 3 shows a domain name system used in an implementation of the present invention;

FIG. 4 shows front-end components of FIG. 2 in greater detail;

FIG. 5 shows back-end components of FIG. 2 in greater detail; and

FIG. 6 and FIG. 7 show conceptual block diagrams of the system of FIG. 2 in alternative contexts.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

The present invention is illustrated and described in terms of a distributed computing environment such as an enterprise computing system using public communication channels such as the Internet. However, an important feature of the present invention is that it is readily scaled upwardly and downwardly to meet the needs of a particular application. Accordingly, unless specified to the contrary, the present invention is applicable to significantly larger, more complex network environments as well as small network environments such as conventional LAN systems.

The present invention is particularly useful in applications where there is an large amount of data

communicated between web servers and web clients (i.e., browser software) or where timeliness (e.g., low latency transport) is important. For example, real-time stock quotes, multi-player games, multi-tiered service to ASP
5 (Application Service Provider) software distribution models benefit from the improvements provided by the present invention. Although the present invention will be described in terms of particular applications, these examples are provided to enhance understanding and are not
10 a limitation of the essential teachings of the present invention. For example, the present invention is readily extended to wireless environments in which the client is a wireless phone or personal digital assistant accessing network resources through a wireless access provider
15 (WAP). In such environments, web-based protocols (i.e., HTTP) may not be used at all. Moreover, while the particular embodiments describe a system that transports web-based traffic using HTTP protocol, it is readily modified to support any public or proprietary protocols
20 including file transfer protocol (FTP), network news protocol (NNTP), mail protocols such as SMTP, voice over Internet protocol (VoIP), and the like.

The current experience of using an overloaded web site that is one of frustration and constant and
25 unexplained errors. In terms of conventional phone-order shopping, this is akin to a customer phoning a favorite mail order company only to receive a constant busy signal or non-stop ringing. Under the present invention, web site operators have a mechanism to communicate with their
30 customers even if they cannot handle their web query immediately. The present invention is able to present the user with a pre-defined set of web pages which can apologize for the delay, instruct the user on what to do in case of an emergency, advertise additional options or

new product offerings, and provide other call-center type features including queue management.

For purposes of this document, a web server is a computer running server software coupled to the World Wide Web (i.e., "the web") that delivers or serves web pages. The web server has a unique IP address and accepts connections in order to service requests by sending back responses. A web server differs from a proxy server or a gateway server in that a web server has resident a set of resources (i.e., software programs, data storage capacity, and/or hardware) that enable it to execute programs to provide an extensible range of functionality such as generating web pages, accessing remote network resources, analyzing contents of packets, reformatting request/response traffic and the like using the resident resources. In contrast, a proxy simply forwards request/response traffic on behalf of a client to resources that reside elsewhere, or obtains resources from a local cache if implemented. A web server in accordance with the present invention may reference external resources of the same or different type as the services requested by a user, and reformat and augment what is provided by the external resources in its response to the user. Commercially available web server software includes Microsoft Internet Information Server (IIS), Netscape Netsite, Apache, among others. Alternatively, a web site may be implemented with custom or semi-custom software that supports HTTP traffic.

FIG. 1 shows an exemplary computing environment 100 in which the present invention may be implemented. Environment 100 includes a plurality of local networks such as Ethernet network 102, FDDI network 103 and Token Ring network 104. Essentially, a number of computing

devices and groups of devices are interconnected through a network 101. For example, local networks 102, 103 and 104 are each coupled to network 101 through routers 109. LANs 102, 103 and 104 may be implemented using any available topology and may implement one or more server technologies including, for example a UNIX, Novell, or Windows NT networks, and/or peer-to-peer type networking. Each network will include distributed storage implemented in each device and typically includes some mass storage device coupled to or managed by a server computer. Network 101 comprises, for example, a public network, such as the Internet, or another network mechanism, such as a fibre channel fabric or conventional WAN technologies.

Local networks 102, 103 and 104 include one or more network appliances 107. One or more network appliances 107 may be configured as an application and/or file server. Each local network 102, 103 and 104 may include a number of shared devices (not shown) such as printers, file servers, mass storage and the like. Similarly, devices 111 may be shared through network 101 to provide application and file services, directory services, printing, storage, and the like. Routers 109 which exist throughout network 101 as well as at the edge of network 101 as shown in FIG. 1, provide a physical connection between the various devices through network 101. Routers 109 may implement desired access and security protocols to manage access through network 101.

Network appliances 107 may also couple to network 101 through public switched telephone network (PSTN) 108 using copper or wireless connection technology. In a typical environment, an Internet service provider 106 supports a connection to network 101 as well as PSTN 108 connections to network appliances 107.

Network appliances 107 may be implemented as any kind of network appliance having sufficient computational function to execute software needed to establish and use a connection to network 101. Network appliances 107 may
5 comprise workstation and personal computer hardware executing commercial operating systems such as Unix variants, Microsoft Windows, Macintosh OS, and the like. At the same time, some appliances 107 comprise portable or handheld devices such as personal digital assistants and
10 cell phones executing operating system software such as PalmOS, WindowsCE, and the like. Moreover, the present invention is readily extended to network devices such as office equipment, vehicles, and personal communicators that occasionally connect through network 101.

Each of the devices shown in FIG. 1 may include memory, mass storage, and a degree of data processing capability sufficient to manage their connection to network 101. The computer program devices in accordance with the present invention are implemented in the memory
15 of the various devices shown in FIG. 1 and enabled by the data processing capability of the devices shown in FIG. 1. In addition to local memory and storage associated with each device, it is often desirable to provide one or more locations of shared storage such as disk farm (not shown)
20 that provides mass storage capacity beyond what an individual device can efficiently use and manage. Selected components of the present invention may be stored in or implemented in shared mass storage.

The present invention operates in a manner akin to a private network 200 implemented within the Internet
30 infrastructure. Private network 200 expedites and prioritizes communications between a client 205 and a web site 210. In the specific examples herein client 205

comprises a network-enabled graphical user interface such as a web browser. However, the present invention is readily extended to client software other than conventional web browser software. Any client application

5 that can access a standard or proprietary user level protocol for network access is a suitable equivalent. Examples include client applications for file transfer protocol (FTP) services, voice over Internet protocol (VoIP) services, network news protocol (NNTP) services,

10 multi-purpose Internet mail extensions (MIME) services, post office protocol (POP) services, simple mail transfer protocol (SMTP) services, as well as Telnet services. In addition to network protocols, the client application may access a network application such as a database management

15 system (DBMS) in which case the client application generates query language (e.g., structured query language or "SQL") messages. In wireless appliances, a client application may communicate via wireless application protocol (WAP) or the like.

20 For convenience, the term "web site" is used interchangeably with "web server" in the description herein, although it should be understood that a web site comprises a collection of content, programs and processes implemented on one or more web servers. A web site is

25 owned by the content provider such as an e-commerce vendor, whereas a web server refers to set of programs running on one or more machines coupled to an Internet node. The web site 210 may be hosted on the site owner's own web server, or hosted on a web server owned by a third

30 party. A web hosting center is an entity that implements one or more web sites on one or more web servers using shared hardware and software resources across the multiple web sites. In a typical web infrastructure, there are many web browsers, each of which has a TCP connection to

the web server in which a particular web site is implemented. The present invention adds two components to the infrastructure: a front-end 201 and back-end 203. Front-end 201 and back-end 203 are coupled by a managed data communication link 202 that forms, in essence, a private network.

Front-end mechanism 201 serves as an access point for client-side communications. Front-end 201 implements a gateway that functions as a proxy for the web server(s) implementing web site 210 (i.e., from the perspective of client 205, front-end 201 appears to be the web site 210). Front-end 201 comprises, for example, a computer that sits "close" to clients 205. By "close", it is meant that the average latency associated with a connection between a client 205 and a front-end 201 is less than the average latency associated with a connection between a client 205 and a web site 210. Desirably, front-end computers have as fast a connection as possible to the clients 205. For example, the fastest available connection may be implemented in point of presence (POP) of an Internet service provider (ISP) 106 used by a particular client 205. However, the placement of the front-ends 201 can limit the number of browsers that can use them. Because of this, in some applications it is more practical to place one front-end computer in such a way that several POPs can connect to it. Greater distance between front-end 201 and clients 205 may be desirable in some applications as this distance will allow for selection amongst a greater number front-ends 201 and thereby provide significantly different routes to a particular back-end 203. This may offer benefits when particular routes and/or front-ends become congested or otherwise unavailable.

Transport mechanism 202 is implemented by cooperative actions of the front-end 201 and back-end 203. Back-end 203 processes and directs data communication to and from web site 210. Transport mechanism 202 communicates data packets using a proprietary protocol over the public Internet infrastructure in the particular example. Hence, the present invention does not require heavy infrastructure investments and automatically benefits from improvements implemented in the general purpose network 101. Unlike the general purpose Internet, front-end 201 and back-end 203 are programmably assigned to serve accesses to a particular web site 210 at any given time.

It is contemplated that any number of front-end and back-end mechanisms may be implemented cooperatively to support the desired level of service required by the web site owner. The present invention implements a many-to-many mapping of front-ends to back-ends. Because the front-end to back-end mappings can be dynamically changed, a fixed hardware infrastructure can be logically reconfigured to map more or fewer front-ends to more or fewer back-ends and web sites or servers as needed.

Front-end 201 together with back-end 203 function to reduce traffic across the TMP link 202 and to improve response time for selected browsers. Traffic across the TMP link 202 is reduced by compressing data and serving browser requests from cache for fast retrieval. Also, the blending of request datagrams results in fewer request:acknowledge pairs across the TMP link 202 as compared to the number required to send the packets individually between front-end 201 and back-end 203. This action reduces the overhead associated with transporting a given amount of data, although conventional request:acknowledge traffic is still performed on the

links coupling the front-end 201 to client 205 and back-end 203 to a web server. Moreover, resend traffic is significantly reduced further reducing the traffic. Response time is further improved for select privileged users and for specially marked resources by determining the priority for each HTTP transmission.

In one embodiment, front-end 201 and back-end 203 are closely coupled to the Internet backbone. This means they have high bandwidth connections, can expect fewer hops, and have more predictable packet transit time than could be expected from a general-purpose connection. Although it is preferable to have low latency connections between front-ends 201 and back-ends 203, a particular strength of the present invention is its ability to deal with latency by enabling efficient transport and traffic prioritization. Hence, in other embodiments front-end 201 and/or back-end 203 may be located farther from the Internet backbone and closer to clients 205 and/or web servers 210. Such an implementation reduces the number of hops required to reach a front-end 201 while increasing the number of hops within the TMP link 202 thereby yielding control over more of the transport path to the management mechanisms of the present invention.

Clients 205 no longer conduct all data transactions directly with the web server 210. Instead, clients 205 conduct some and preferably a majority of transactions with front-ends 201, which simulate the functions of web server 210. Client data is then sent, using TMP link 202, to the back-end 203 and then to the web server 210. Running multiple clients 205 over one large connection provides several advantages:

- Since all client data is mixed, each client can be assigned a priority. Higher priority clients, or

clients requesting higher priority data, can be given preferential access to network resources so they receive access to the channel sooner while ensuring low-priority clients receive sufficient service to meet their needs.

- The large connection between a front-end 201 and back-end 203 can be permanently maintained, shortening the many TCP/IP connection sequences normally required for many clients connecting and disconnecting.

Using a proprietary protocol allows the use of more effective techniques to improve data throughput and makes better use of existing bandwidth during periods when the network is congested.

A particular advantage of the architecture shown in FIG. 2 is that it is readily scaled. Any number of client machines 205 may be supported. In a similar manner, a web site owner may choose to implement a site using multiple web servers 210 that are co-located or distributed throughout network 101. To avoid congestion, additional front-ends 201 may be implemented or assigned to particular web sites. Each front-end 201 is dynamically re-configurable by updating address parameters to serve particular web sites. Client traffic is dynamically directed to available front-ends 201 to provide load balancing. Hence, when quality of service drops because of a large number of client accesses, an additional front-end 201 can be assigned to the web site and subsequent client requests directed to the newly assigned front-end 201 to distribute traffic across a broader base.

In the particular examples, this is implemented by a front-end manager component 207 that communicates with

multiple front-ends 201 to provide administrative and configuration information to front-ends 201. Each front-end 201 includes data structures for storing the configuration information, including information
5 identifying the IP addresses of web servers 210 to which they are currently assigned. Other administrative and configuration information stored in front-end 201 may include information for prioritizing data from and to particular clients, quality of service information, and
10 the like.

Similarly, additional back-ends 203 can be assigned to a web site to handle increased traffic. Back-end manager component 209 couples to one or more back-ends 203 to provide centralized administration and configuration
15 service. Back-ends 203 include data structures to hold current configuration state, quality of service information and the like. In the particular examples front-end manager 207 and back-end manager 209 serve multiple web sites 210 and so are able to manipulate the
20 number of front-ends and back-ends assigned to each web site 210 by updating this configuration information. When the congestion for the site subsides, the front-end 201 and back-end 203 can be reassigned to other, busier web sites. These and similar modifications are equivalent to
25 the specific examples illustrated herein.

In the case of web-based environments, front-end 201 is implemented using custom or off-the-shelf web server software. Front-end 201 is readily extended to support other, non-web-based protocols, however, and may support
30 multiple protocols for varieties of client traffic. Front-end 201 processes the data traffic it receives, regardless of the protocol of that traffic, to a form suitable for transport by TMP 202 to a back-end 203.

Hence, most of the functionality implemented by front-end 201 is independent of the protocol or format of the data received from a client 205. Hence, although the discussion of the exemplary embodiments herein relates primarily to front-end 201 implemented as a web server, it should be noted that, unless specified to the contrary, web-based traffic management and protocols are merely examples and not a limitation of the present invention.

As shown in FIG. 2, in accordance with the present invention a web site is implemented using an originating web server 210 operating cooperatively with the web server of front-end 201. More generally, any network service (e.g., FTP, VoIP, NNTP, MIME, SMTP, Telnet, DBMS) can be implemented using a combination of an originating server working cooperatively with a front-end 201 configured to provide a suitable interface (e.g., FTP, VoIP, NNTP, MIME, SMTP, Telnet, DBMS, WAP) for the desired service. In contrast to a simple front-end cache or proxy software, implementing a server in front-end 201 enables portions of the web site (or other network service) to actually be implemented in and served from both locations. The actual web pages or service being delivered comprises a composite of the portions generated at each server. Significantly, however, the web server in front-end 201 is close to the browser in a client 205 whereas the originating web server is close to all resources available at the web hosting center at which web site 210 is implemented. In essence the web site 210 is implemented by a tiered set of web servers comprising a front-end server 201 standing in front of an originating web server.

This difference enables the web site or other network service to be implemented so as to take advantage of the unique topological position each entity has with respect

to the client 205. By way of a particular example, assume an environment in which the front-end server 201 is located at the location of an ISP used by a particular set of clients 205. In such an environment, clients 205 can
5 access the front-end server 205 without actually traversing the network 101.

In order for a client 205 to obtain service from a front-end 201, it must first be directed to a front-end 201 that can provide the desired service. Preferably,
10 client 205 does not need to be aware of the location of front-end 201, and initiates all transactions as if it were contacting the originating server 210. FIG. 3 illustrates a domain name server (DNS) redirection mechanism that illustrates how a client 205 is connected
15 to a front-end 201. The DNS systems is defined in a variety of Internet Engineering Task Force (IETF) documents such as RFC0883, RFC 1034 and RFC 1035 which are incorporated by reference herein. In a typical environment, a client 205 executes a browser 301, TCP/IP
20 stack 303, and a resolver 305. For reasons of performance and packaging, browser 301, TCP/IP stack 303 and resolver 305 are often grouped together as routines within a single software product.

Browser 301 functions as a graphical user interface
25 to implement user input/output (I/O) through monitor 311 and associated keyboard, mouse, or other user input device (not shown). Browser 301 is usually used as an interface for web-based applications, but may also be used as an interface for other applications such as email and network
30 news, as well as special-purpose applications such as database access, telephony, and the like. Alternatively, a special-purpose user interface may be substituted for

the more general purpose browser 301 to handle a particular application.

TCP/IP stack 303 communicates with browser 301 to convert data between formats suitable for browser 301 and IP format suitable for Internet traffic. TCP/IP stack also implements a TCP protocol that manages transmission of packets between client 205 and an Internet service provider (ISP) or equivalent access point. Internet protocol (IP) requires that each data packet include, among other things, an IP address identifying a destination node. In current implementations the IP address comprises a 32-bit value that identifies a particular Internet node. Non-IP networks have similar node addressing mechanisms. To provide a more user-friendly addressing system, the Internet implements a system of domain name servers that map alpha-numeric domain names to specific IP addresses. This system enables a name space that is more consistent reference between nodes on the Internet and avoids the need for users to know network identifiers, addresses, routes and similar information in order to make a connection.

The domain name service is implemented as a distributed database managed by domain name servers (DNSs) 307 such as DNS_A, DNS_B and DNS_C shown in FIG. 3. Each DNS relies on <domain name:IP> address mapping data stored in master files scattered through the hosts that use the domain system. These master files are updated by local system administrators. Master files typically comprise text files that are read by a local name server, and hence become available through the name servers 307 to users of the domain system.

The user programs (e.g., clients 205) access name servers through standard programs such as resolver 305.

Resolver 305 includes an address of a DNS 307 that serves as a primary name server. When presented with a reference to a domain name (e.g., <http://www.circadence.com>), resolver 305 sends a request to the primary DNS (e.g., DNS_A in FIG. 3). The primary DNS 307 returns either the IP address mapped to that domain name, a reference to another DNS 307 which has the mapping information (e.g., DNS_B in FIG. 3), or a partial IP address together with a reference to another DNS that has more IP address information. Any number of DNS-to-DNS references may be required to completely determine the IP address mapping.

In this manner, the resolver 305 becomes aware of the IP address mapping which is supplied to TCP/IP component 303. Client 205 may cache the IP address mapping for future use. TCP/IP component 303 uses the mapping to supply the correct IP address in packets directed to a particular domain name so that reference to the DNS system need only occur once.

In accordance with the present invention, at least one DNS server 307 is owned and controlled by system components of the present invention. When a user accesses a network resource (e.g., a web site), browser 301 contacts the public DNS system to resolve the requested domain name into its related IP address in a conventional manner. In a first embodiment, the public DNS performs a conventional DNS resolution directing the browser to an originating server 210 and server 210 performs a redirection of the browser to the system owned DNS server (i.e., DNC_C in FIG. 3). In a second embodiment, domain:address mappings within the DNS system are modified such that resolution of the of the originating server's domain automatically return the address of the system-owned DNS server (DNS_C). Once a browser is redirected to

the system-owned DNS server, it begins a process of further redirecting the browser 301 to the best available front-end 201.

Unlike a conventional DNS server, however, the system-owned DNS_C in FIG. 3 receives domain:address mapping information from a redirector component 309. Redirector 309 is in communication with front-end manager 207 and back-end manager 209 to obtain information on current front-end and back-end assignments to a particular server 210. A conventional DNS is intended to be updated infrequently by reference to its associated master file. In contrast, the master file associated with DNS_C is dynamically updated by redirector 309 to reflect current assignment of front-end 201 and back-end 203. In operation, a reference to web server 210 (e.g., <http://www.circadence.com>) may result in an IP address returned from DNS_C that points to any selected front-end 201 that is currently assigned to web site 210. Likewise, web site 210 may identify a currently assigned back-end 203 by direct or indirect reference to DNS_C.

Front-end 201 typically receives information directly from front-end manager 207 about the address of currently assigned back-ends 203. Similarly, back-end 203 is aware of the address of a front-end 201 associated with each data packet. Hence, reference to the domain system is not required to map a front-end 201 to its appropriate back-end 203.

FIG. 4 illustrates principle functional components of an exemplary front-end 201 in greater detail. Primary functions of the front-end 201 include serving as a proxy for web server 210 from the perspective of client 205, and translating transmission control protocol (TCP) packets from client 205 into transport morphing protocol™ (TMP™)

packets used in the system in accordance with the present invention. Transport morphing protocol and TMP are trademarks or registered trademarks of Circadence corporation in the United States and other countries. It is contemplated that the various functions described in reference to the specific examples may be implemented using a variety of data structures and programs operating at any location in a distributed network. For example, a front-end 201 may be operated on a network appliance 107 or server within a particular network 102, 103, or 104 shown in FIG. 1. The present invention is readily adapted to any application where multiple clients are coupling to a centralized resource. Moreover, other transport protocols may be used, including proprietary transport protocols.

TCP component 401 includes devices for implementing physical connection layer and Internet protocol (IP) layer functionality. Current IP standards are described in IETF documents RFC0791, RFC0950, RFC0919, RFC0922, RFC792, RFC1112 that are incorporated by reference herein. For ease of description and understanding, these mechanisms are not described in great detail herein. Where protocols other than TCP/IP are used to couple to a client 205, TCP component 401 is replaced or augmented with an appropriate network protocol process.

TCP component 401 communicates TCP packets with one or more clients 205. Received packets are coupled to parser 402 where the Internet protocol (or equivalent) information is extracted. TCP is described in IETF RFC0793 which is incorporated herein by reference. Each TCP packet includes header information that indicates addressing and control variables, and a payload portion that holds the user-level data being transported by the

TCP packet. The user-level data in the payload portion typically comprises a user-level network protocol datagram.

Parser 402 analyzes the payload portion of the TCP
5 packet. In the examples herein, HTTP is employed as the user-level protocol because of its widespread use and the advantage that currently available browser software is able to readily use the HTTP protocol. In this case, parser 402 comprises an HTTP parser. More generally,
10 parser 402 can be implemented as any parser-type logic implemented in hardware or software for interpreting the contents of the payload portion. Parser 402 may implement file transfer protocol (FTP), mail protocols such as simple mail transport protocol (SMTP), structured query
15 language (SQL) and the like. Any user-level protocol, including proprietary protocols, may be implemented within the present invention using appropriate modification of parser 402.

To improve performance, front-end 201 optionally
20 includes a caching mechanism 403. Cache 403 may be implemented as a passive cache that stores frequently and/or recently accessed web pages or as an active cache that stores network resources that are anticipated to be accessed. In non-web applications, cache 403 may be used
25 to store any form of data representing database contents, files, program code, and other information. Upon receipt of a TCP packet, HTTP parser 402 determines if the packet is making a request for data within cache 403. If the request can be satisfied from cache 403, the data is
30 supplied directly without reference to web server 210 (i.e., a cache hit). Cache 403 implements any of a range of management functions for maintaining fresh content. For example, cache 403 may invalidate portions of the

cached content after an expiration period specified with the cached data or by web sever 210. Also, cache 403 may proactively update the cache contents even before a request is received for particularly important or frequently used data from web server 210. Cache 403 evicts information using any desired algorithm such as least recently used, least frequently used, first in/first out, or random eviction. When the requested data is not within cache 403, a request is processed to web server 210, and the returned data may be stored in cache 403.

Several types of packets will cause parser 404 to forward a request towards web server 210. For example, a request for data that is not within cache 403 (or if optional cache 403 is not implemented) will require a reference to web server 210. Some packets will comprise data that must be supplied to web server 210 (e.g., customer credit information, form data and the like). In these instances, HTTP parser 402 couples to data blender 404.

Optionally, front-end 201 implements security processes, compression processes, encryption processes and the like to condition the received data for improved transport performance and/or provide additional functionality. These processes may be implemented within any of the functional components (e.g., data blender 404) or implemented as separate functional components within front-end 201. Also, parser 402 may implement a prioritization program to identify packets that should be given higher priority service. A prioritization program requires only that parser 402 include a data structure associating particular clients 205, particular TCP packet types, contents and the like with a prioritization value. Based on the prioritization value, parser 402 may

selectively implement such features as caching, encryption, security, compression and the like to improve performance and/or functionality. The prioritization value is provided by the owners of web site 210, for example, and may be dynamically altered, statically set, or updated from time to time to meet the needs of a particular application.

Blender 404 slices and/or coalesces the data portions of the received packets into a more desirable "TMP units" that are sized for transport through the TMP mechanism 202. The data portion of TCP packets may range in size depending on client 205 and any intervening links coupling client 205 to TCP component 401. Moreover, where compression is applied, the compressed data will vary in size depending on the compressibility of the data. Data blender 404 receives information from front-end manager 207 that enables selection of a preferable TMP packet size. Alternatively, a fixed TMP packet size can be set that yields desirable performance across TMP mechanism 202. Data blender 404 also marks the TMP units so that they can be re-assembled at the receiving end.

Data blender 404 also serves as a buffer for storing packets from all clients 205 that are associated with front-end 201. Blender 404 mixes data packets coming into front-end 201 into a cohesive stream of TMP packets sent to back-end 203 over TMP link 202. In creating a TMP packet, blender 404 is able to pick and choose amongst the available data packets so as to prioritize some data packets over others.

In an exemplary implementation, a "TMP connection" comprises a plurality of "TCP connection buffers", logically arranged in multiple "rings". Each TCP socket maintained between the front-end 201 and a client 205

corresponds to a TCP connection buffer. When a TCP connection buffer is created, it is assigned a priority. For purposes of the present invention, any algorithm or criteria may be used to assign a priority. Each priority ring is associated with a number of TCP connection buffers having similar priority. In a specific example, five priority levels are defined corresponding to five priority rings. Each TMP ring is characterized by the number of connection buffers it holds (nSockets), the number of connection buffers it holds that have data waiting to be sent (nReady) and the total number of bytes of data in all the connection buffers that it holds (nBytes).

When composing TMP data packets, the blender goes into a loop comprising the steps:

1) Determine the number of bytes available to be sent from each ring (nBytes), and the number of TCP connections that are ready to send (nReady)

2) Determine how many bytes should be sent from each ring. This is based on a weight parameter for each priority. The weight can be thought of as the number of bytes that should be sent at each priority this time through the loop.

3) The nSend value computed in the previous step reflects the weighted proportion that each ring will have in a blended TMP packet, but the values of nSend do not reflect how many bytes need to be selected to actually empty most or all of the data waiting to be sent a single round. To do this, the nSend value is normalized to the ring having the most data waiting (e.g., nBytes = nSendNorm). This involves a calculation of a factor: $S = nBytes / (Weight * nReady)$ for the ring with the greatest nReady. Then, for each ring, calculate $nReady * S * Weight$ to

get the normalized value (nSendNorm) for each priority ring.

4) Send sub-packets from the different rings. This is done by taking a sub-packet from the highest priority ring and adding it to a TMP packet, then adding a sub-packet from each of the top two queues, then the top three, and so on.

5) Within each ring, sub-packets are added round robin. When a sub-packet is added from a TCP connection buffer the ring is rotated so the next sub-packet the ring adds will come from a different TCP connection buffer. Each sub-packet can be up to 512 bytes in a particular example. If the connection buffer has less than 512 bytes waiting, the data available is added to the TMP packet.

6) When a full TMP packet (roughly 1.5 kB in a particular example) is built, it is sent. This can have three or more sub packets, depending on their size. The TMP packet will also be sent when there is no more data ready.

20 TMP mechanism 405 implements the TMP protocol in accordance with the present invention. TMP is a TCP-like protocol adapted to improve performance for multiple channels operating over a single connection. Front-end TMP mechanism 405 and a corresponding back-end TMP
25 mechanism 505 shown in FIG. 5 are computer processes that implement the end points of TMP link 202. The TMP mechanism in accordance with the present invention creates and maintains a stable connection between two processes for high-speed, reliable, adaptable communication.

30 TMP is not merely a substitute for the standard TCP environment. TMP is designed to perform particularly well

in heterogeneous environments such as the Internet. TMP connections are made less often than TCP connections. Once a TMP connection is made, it remains up unless there is some kind of direct intervention by an administrator or there is some form of connection breaking network error. This reduces overhead associated with setting up, maintaining and tearing down connections normally associated with TCP.

Another feature of TMP is its ability to channel numerous TCP connections through a single TMP connection 202. The environment in which TMP resides allows multiple TCP connections to occur at one end of the system. These TCP connections are then mapped into a single TMP connection. The TMP connection is then broken down at the other end of the TMP connection 202 in order to traffic the TCP connections to their appropriate destinations. TMP includes mechanisms to ensure that each TMP connection gets enough of the available bandwidth to accommodate the multiple TCP connections that it is carrying.

Another advantage of TMP as compared to traditional protocols is the amount of information about the quality of the connection that a TMP connection conveys from one end to the other of a TMP connection 202. As often happens in a network environment, each end has a great deal of information about the characteristics of the connection in one direction, but not the other. By knowing about the connection as a whole, TMP can better take advantage of the available bandwidth.

In contrast with conventional TCP mechanisms, the behavior implemented by TMP mechanism 405 is constantly changing. Because TMP obtains bandwidth to host a variable number of TCP connections and because TMP is responsive information about the variable status of the

network, the behavior of TMP is preferably continuously variable. One of the primary functions of TMP is being able to act as a conduit for multiple TCP connections. As such, a single TMP connection cannot behave in the same manner as a single TCP connection. For example, imagine that a TMP connection is carrying 100 TCP connections. At this time, it loses one packet (from any one of the TCP connections) and quickly cuts its window size in half (as specified for TCP). This is a performance reduction on 100 connections instead of just on the one that lost the packet.

Each TCP connection that is passed through the TMP connection must get a fair share of the bandwidth, and should not be easily squeezed out. To allow this to happen, every TMP becomes more aggressive in claiming bandwidth as it accelerates. Like TCP, the bandwidth available to a particular TMP connection is measured by its window size (i.e., the number of outstanding TCP packets that have not yet been acknowledged). Bandwidth is increased by increasing the window size, and relinquished by reducing the window size. Up to protocol specified limits, each time a packet is successfully delivered and acknowledged, the window size is increased until the window size reaches a protocol specified maximum. When a packet is dropped (e.g., no acknowledge received or a resend packet response is received), the bandwidth is decreased by backing off the window size. TMP also ensures that it becomes more and more resistant to backing off (as compared to TCP) with each new TCP connection that it hosts. A TMP should not go down to a window size of less than the number of TCP connections that it is hosting.

In a particular implementation, every time a TCP connection is added to (or removed from) what is being passed through the TMP connection, the TMP connection behavior is altered. It is this adaptation that ensures successful connections using TMP. Through the use of the adaptive algorithms discussed above, TMP is able to adapt the amount of bandwidth that it uses. When a new TCP connection is added to the TMP connection, the TMP connection becomes more aggressive. When a TCP connection is removed from the TMP connection, the TMP connection becomes less aggressive.

5
10
15
20
25
TMP connection 202 provides improved performance in its environment as compared to conventional TCP channels, but it is recognized that TMP connection 202 resides on the open, shared Internet in the preferred implementations. Hence, TMP must live together with many protocols and share the connection efficiently in order to allow the other transport mechanisms fair access to the shared communication bandwidth. Since TMP takes only the amount of bandwidth that is appropriate for the number of TCP connections that it is hosting (and since it monitors the connection and controls the number of packets that it puts on the line), TMP will exist cooperatively with TCP traffic. Furthermore, since TMP does a better job at connection monitoring than TCP and TMP is better suited to throughput and bandwidth management than TCP.

Also shown in FIG. 4 are data filter component 406 and HTTP reassemble component 407 that process incoming (with respect to client 205) data. TMP mechanism 405 receives TMP packets from TMP connection 202 and extracts the TMP data units. Using the appended sequencing information, the extracted data units are reassembled into HTTP data packet information by HTTP reassembler 407.

Data filter component 406 may also implement data decompression where appropriate, decryption, and handle caching when the returning data is of a cacheable type.

FIG. 5 illustrates principle functional components of an exemplary back-end 203 in greater detail. Primary functions of the back-end 203 include translating transmission control protocol (TCP) packets from web server 210 into TMP packets as well as translating TMP packets received from a front-end 201 into the one or more corresponding TCP packets to be send to server 210.

10 TMP unit 505 receives TMP packets from TMP connection 202 and passes them to HTTP reassemble unit 507 where they are reassembled into the corresponding TCP packets. Data filter 506 may implement other functionality such as decompression, decryption, and the like to meet the needs of a particular application. The reassembled data is forwarded to TCP component 501 for communication with web server 210.

20 TCP data generated by the web server process are transmitted to TCP component 501 and forwarded to HTTP parse mechanism 502. Parser 502 operates in a manner analogous to parser 402 shown in FIG. 4 to extract the data portion from the received TCP packets, perform optional compression, encryption and the like, and forward those packets to data blender 504. Data blender 504 operates in a manner akin to data blender 404 shown in FIG. 3 to buffer and prioritize packets in a manner that is efficient for TMP transfer. Priority information is received by, for example, back-end manager 209 based upon criteria established by the web site owner. TMP data is streamed into TMP unit 505 for communication on TMP connection 202.

To improve performance, back-end 203 optionally includes a caching mechanism 503. Cache 503 may be implemented as a passive cache that stores frequently and/or recently accessed resources from server 210 or as an active cache that stores network resources that are anticipated to be accessed. In non-web applications, cache 503 may be used to store any form of data representing database contents, files, program code, and other information. Upon receipt of a TCP packet, HTTP reassemble unit 507 determines if the packet is making a request for data within cache 503. If the request can be satisfied from cache 503, the data is supplied directly without reference to web server 210. Cache 503 implements any of a range of management functions for maintaining fresh content. For example, cache 503 may invalidate portions of the cached content after an expiration period specified with the cached data or by web sever 210. Also, cache 503 may proactively update the cache contents even before a request is received for particularly important or frequently used data from web server 210. Cache 503 evicts information using any desired algorithm such as least recently used, least frequently used, first in/first out, or random eviction. When the requested data is not within cache 503, a reassembled HTTP request is processed to web server 210, and the returned data may be stored in cache 503.

Returning again to FIG. 2, in a particular example front-end server 201 comprises a web server or hosting center located in a first geographic region whereas web site 210 comprises a web server or hosting center located in a second geographic region. The Internet is organized as a collection of subnets often arranged around roughly geographic boundaries. Latency and quality of service within a subnet is often better than between subnets. A

client 205 located in the first geographic region principally connects to front-end 201, but occasionally requests content that is available only on web site 210. In conventional implementations, a web server 201
5 redirects client 205 to web server 210.

In accordance with the present invention, front-end 201 retains control of the connection with client 205 and selectively obtains content from web server 210 or serves substitute content. Front-end server 201 serves as a
10 first tier web server which is backed by a second tier web server 210. This action protects client 205 from the vagaries of the subnet-to-subnet latency and quality of service and avoids undesirable handoff of a client 205 to a logically distant web server 210.

FIG. 6 illustrates a conceptual block diagram of an exemplary implementation of the system shown in FIG. 2 in an alternative context. In the example of FIG. 6, front-end 201 is implemented as a front-end server 601 operating at, for example, an ISP 602. ISP 602 supports modem,
20 digital subscriber line (DSL), ISDN, leased line, or other communication ports for communicating with one or more clients 605. ISP 602 serves as a bridge to couple client connections to IP connections with network 101. ISP 602 may be considered to be outside of network 101 in that
25 quality of service between client 605 and ISP 602 is not dependent on congestion or equipment failure or other factors affecting quality of service within network 101.

In operation, client 605 generates an HTTP request specifying web site 610 in the URL. In the manner
30 described hereinbefore, the client request is redirected to front-end 601. Once the redirection is completed, front-end 601 serves web pages embedded in HTTP response packets using both content obtained from web site 610 as

well as content obtained from front-end content database 603. In this manner, front-end server 601 dynamically controls the source of the delivered content. The web page(s) served to client 605 comprise a composite of multiple sources from multiple independent web servers.

Significantly, the content served from content database 603 may differ from the content that would have otherwise been served by web site 610. For example, if web site 610 returns an HTTP 404 "page not found" error page, front-end site 601 may supply a more informative or instructive web page derived from content database 603. Alternatively, front-end site 601 may detect periods of low quality of service or slow response and provide substitute content from content database 603 or elsewhere. In a particular example, web site 610 publishes a load index that can be read by front-end 601 and used to generate a wait page intended to occupy a user of client 605 until content can be obtained directly from web site 610.

In essence this substitute content implements a "soft landing" feature that provides improved behavior in the event of web server outages, slowdowns, and/or traffic congestion. The soft landing behavior may include serving substitute content that indicates an estimated time for server availability, demonstration of product features, or provides a coupon for discounts on future goods and services through the site. The soft landing behavior may also include a manual or automatic redirection to a mirror server. The actual soft landing behavior might change depending on priority value associated with a particular user. For example, a high priority user might be given prioritized access to the origin web site while the low priority user is served substitute content. When a coupon

is served as a soft landing, the coupon value or terms can be varied to discriminate between high and low priority users.

In a particular example, the owner of web site 610 establishes rules for how front-end server 601 handles various conditions. These rules are stored in front-end server 601. Hence, the owner of web site 610 is not losing control over how web site 610 is presented, but instead is gaining control over how presentation occurs during periods where network 101 is unavailable or provides unacceptable quality of service. By placing a web server in front of the origin web site 610, overall user experience as well as efficacy of the web site 610 for the site owner are improved.

Content database 603 may be implemented in whole or part by cache 403 shown in Fig. 4. It is contemplated that a back-end server 203 may be used to provide the alternate content serving features described in reference to Fig. 6 in which case back-end cache 503 provides a content database 603. Moreover, a combination of front-end server 201 and back-end server 203 together can be used. In each alternative, caching can be passive or active. Passive caching will cache network resources as they pass through the server. Active caching can be implemented by causing the server to examine the network resources as it passes through the server to identify caching instructions encoded within the network resource itself. These instructions specify whether the resource is cacheable, expiration times, and the like. These instructions may be provided by the web server 210 itself, or may be provided to a front-end 201 by a back-end 203 according to rules specified by the site owner or other rule making entity.

Fig. 7 shows yet another alternative in which the system of the present invention may be implemented. By way of contrast, Fig. 2 illustrates an implementation using coordinated actions of a front-end server 201, a back-end server 203, and one or more web servers 210. In contrast, Fig. 6 illustrates a single intermediary server 601 that cooperating with a web server 610. Fig. 7, however, illustrates a multi-tiered web server approach in which any number of intermediary servers 701 act cooperatively and dynamically to serve a web site 710 in a distributed fashion. A browser client 705 couples to one, or more than one web server 701 to assert a request for content. The web server 701 may serve content directly, or may obtain content from another web server 701, or from the web site 710 as in the embodiment shown in Fig. 6. Each intermediary web server 701 exercises an option to serve content directly from its own resources, or obtain the content from another resource such as a web server 701.

The implementation shown in Fig. 7 will benefit from management functions that inform each intermediary web server 701 about the content and/or resources available on each other web server 701. This management information enables requests to be directed to web servers 701 that are likely to be able to answer the request. Routing through multiple tiers of web servers 701 will eventually increase the latency in generating a response as compared to a direct reference to web site 710. However, when web servers 701 are coupled by enhanced channels such as the TCP channels described hereinbefore, it is contemplated that a request may involve multiple tiers of web servers 701 and still be served faster, more reliably than possible through direct reference to web site 710. From another perspective, low priority requests may be

preferably be served from a multiple-tier service involving several intermediary servers 701 such that higher priority requests may be served using fewer tiers or by direct reference to web site 710. Hence, the
5 present invention enables a means for regulating the quality of service provided to a given client request by regulating the number of intermediary servers involved in responding to the request.

Although the invention has been described and
10 illustrated with a certain degree of particularity, it is understood that the present disclosure has been made only by way of example, and that numerous changes in the combination and arrangement of parts can be resorted to by those skilled in the art without departing from the spirit
15 and scope of the invention, as hereinafter claimed. For example, while devices supporting HTTP data traffic are used in the examples, the HTTP devices may be replaced or augmented to support other public and proprietary protocols and languages including FTP, NNTP, SMTP, SQL and
20 the like. In such implementations the front-end 201 and/or back end 203 are modified to implement the desired protocol. Moreover, front-end 201 and back-end 203 may support different protocols and languages such that the front-end 201 supports, for example, HTTP traffic with a
25 client and the back-end supports a DBMS protocol such as SQL. Such implementations not only provide the advantages of the present invention, but also enable a client to access a rich set of network resources with minimal client software.